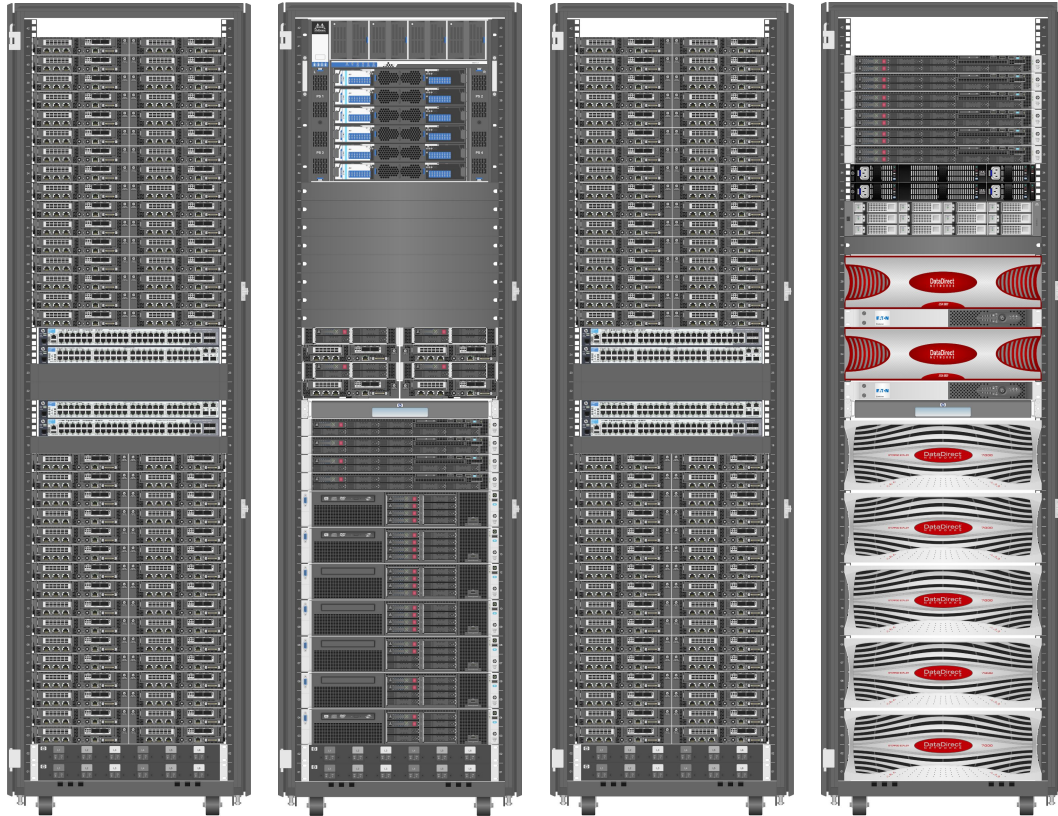


Guía básica para el uso de Leftrararu



Arquitectura Leftrarú



- Infiniband FDR 56Gbps
- Red iLO para adm. hw.
- Red servicio 1 Gbps
- +200 Tb almacenamiento DDN Lustre
- 128 nodos *slims*
- 4 nodos *fats*
- 12 Xeon Phi
- Racks enfriados por agua
- Enfriamiento in-row respaldo
- UPS 120 KVA autonomía: 30 mins.

Nodos de cómputo Leftrarú

2 tipos de nodo de cómputo:

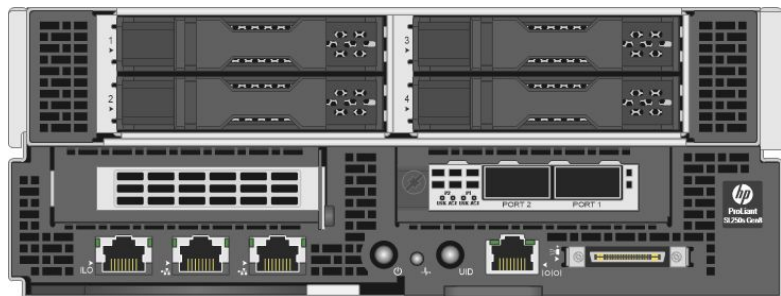
Slim



Nodo Slim

- 2 CPUs Xeon E5-2660 v2, 10 cores c.u.
- 48 GB RAM DIMM DDR3
- HD Interno 300 GB
- InfiniBand FDR 56 Gbps

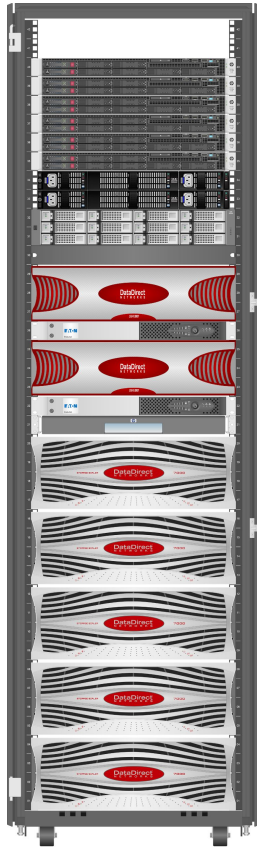
Fat



Nodo Fat

- 2 CPUs Xeon E5-2660 v2, 10 cores c.u.
- 64 GB RAM DIMM DDR3
- HD Interno 300 GB
- InfiniBand FDR 56 Gbps
- 3 coprocesadores Xeon Phi
 - 240 cores c.u.
 - 8 GB RAM c.u.

Almacenamiento Leftraru



DDN EXAScaler

- Almacenamiento paralelo de clase mundial
- Alto rendimiento en operaciones IO
- Tolerante a fallas (alta disponibilidad)
- Interconexión infiniband
- Capacidades Big Data

Características en Leftraru

- Sistema archivos EXAScaler (Lustre)
- +200 TB de almacenamiento
- Almacenamiento metadata separado
- 2 controladoras SFA en H.A.
- 4 nodos OSS conectados a Infiniband
- 2 nodos MDS en H.A.

Introducción a SLURM

Simple Linux Utility for Resource Management

- Administra clusters y ejecución de tareas
- Open source
- Utilizado en el 60% de los supercomputadores del top500
- Versión actual instalada en Leftrarú: 16.05.4



Características de SLURM

- Tres funciones principales:
 - Asignación de recursos (exclusivos y no exclusivos)
 - Framework para iniciar, ejecutar y monitorear trabajos
 - Gestiona tareas manejando una cola de recursos
- Integra una base de datos para reportes históricos
- Puede reservar diferentes recursos: CPU, socket, nodo o incluso por RAM
 - 6 nodos o 120 CPUs en el caso de Leftraru
- Permite a los administradores modificar tareas en ejecución y realizar reservas

Slurm: Primeros pasos

Pruebas desde la consola:

```
[usuario@leftraru1 ~]$ srun -n 2 hostname
```

```
cn028
```

```
cn028
```

```
[usuario@leftraru1 ~]$
```

2 tareas

salida comando



Slurm: Monitorear mis tareas desde consola

Monitorear desde la consola:

```
[usuario@leftraru1 ~]$ squeue
```



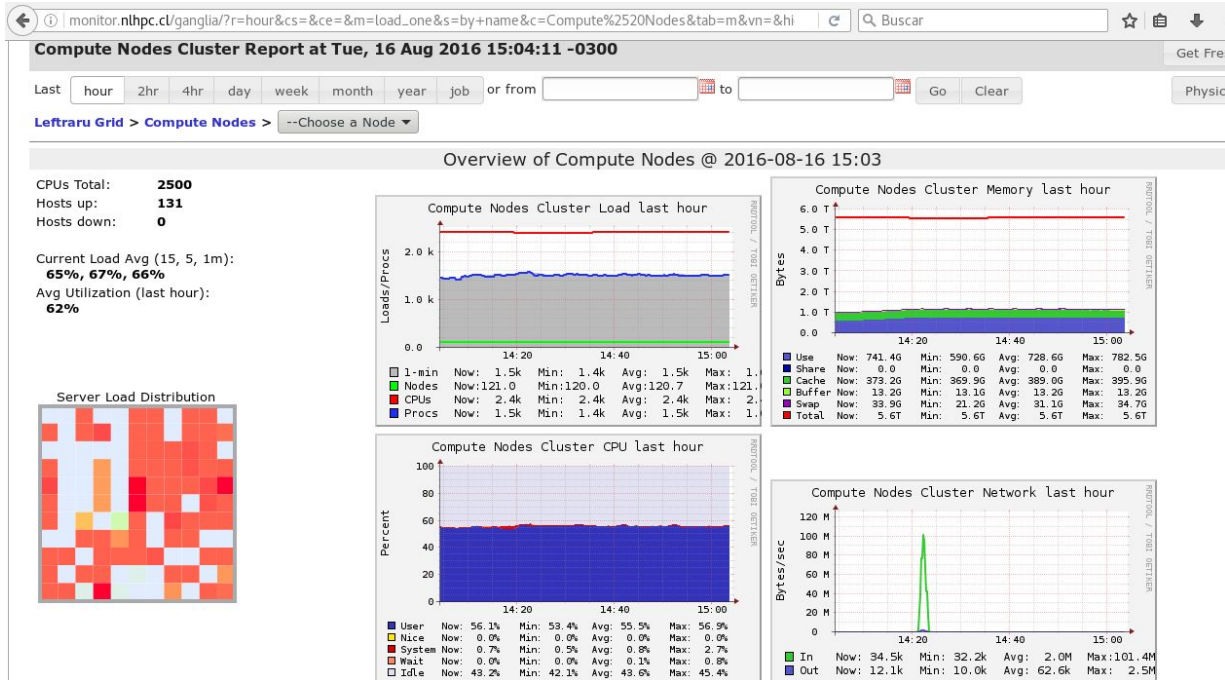
JOBID	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST(REASON)
4400799	slims	example	usuario	R	0:00	1	cn042

Obtener el nombre de el/los nodos y visitar:

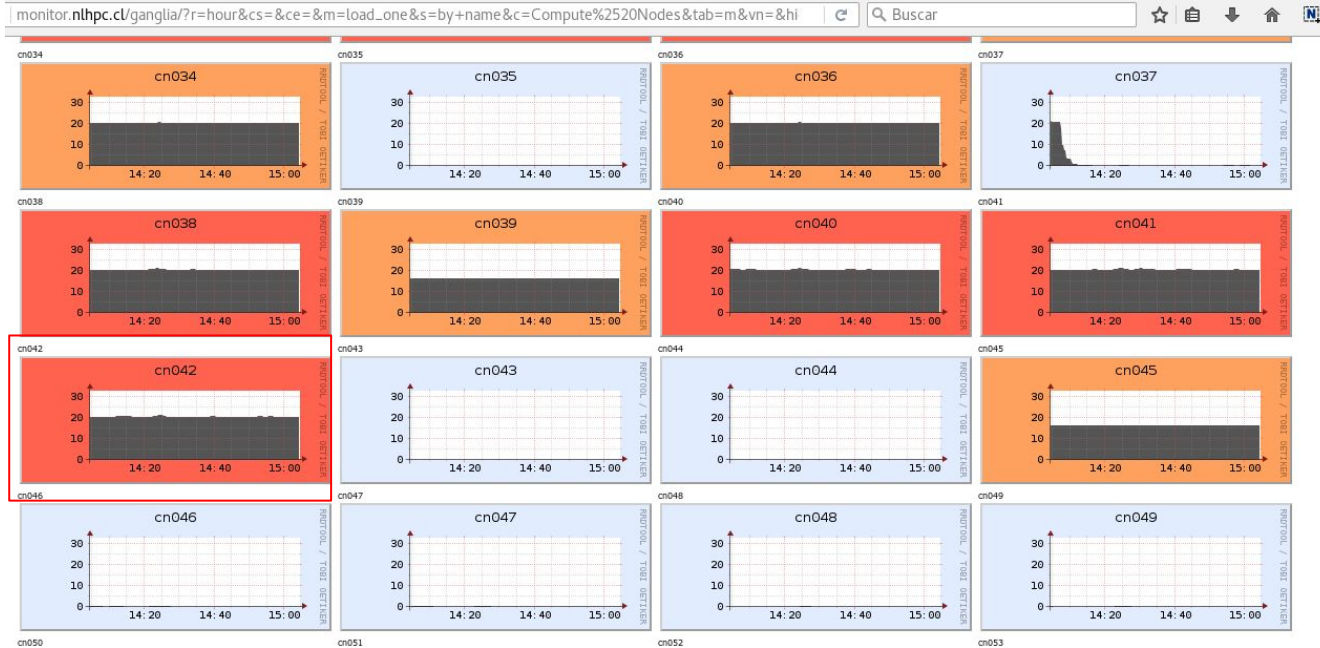
<http://monitor.nlhpc.cl/ganglia>



Slurm: Monitorear mis tareas desde la web



Slurm: Monitorear mis tareas desde la web



Slurm: Monitorear mis tareas desde la web



Monitorear mis tareas

Puede ingresar a través de ssh a un nodo en donde tenga una tarea en ejecución y ejecutar htop

```
1 [|||||||||||||100.0%] 6 [ 0.0%] 11 [ 0.0%] 16 [ 0.0%]
2 [ 0.0%] 7 [ 0.0%] 12 [ 0.0%] 17 [ 0.0%]
3 [|||||||||||||100.0%] 8 [ 0.0%] 13 [ 0.0%] 18 [ 0.0%]
4 [ 0.0%] 9 [ 0.7%] 14 [ 0.0%] 19 [ 0.0%]
5 [ 0.0%] 10 [ 0.0%] 15 [ 0.0%] 20 [ 0.0%]

Mem[||||| 1.56G/62.7G] Tasks: 44, 39 thr: 3 running
Swp[ 0K/62.5G] Load average: 2.00 2.01 2.05
Uptime: 4 days, 00:36:11

  PID USER   PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
18305 nperinet  20   0 19028  3768   892  R 100.0  0.0  2h52:34  /LL RK4.x
18335 nperinet  20   0 17104  1584   892  R 100.0  0.0  2h47:30  /LL RK4.x
18605 root      20   0  121M  2432  1476  R  0.0  0.0  0:00.05  htop
    1 root      20   0  185M  5068  2384  S  0.0  0.0  0:11.06  /usr/lib/systemd/systemd --switched-root --system --de
   669 root      20   0  112M  2000  1564  S  0.0  0.0  0:00.00  /bin/bash
   671 root      20   0 36816  7236  6908  S  0.0  0.0  0:00.78  /usr/lib/systemd/systemd-journald
   726 root      20   0 43808  2428  1268  S  0.0  0.0  0:00.83  /usr/lib/systemd/systemd-udev
  1012 root      16  -4 51188  1620  1236  S  0.0  0.0  0:00.05  /sbin/auditd -n
  1002 root      16  -4 51188  1620  1236  S  0.0  0.0  0:00.25  /sbin/auditd -n
  1206 root      20   0  448M 10956  6968  S  0.0  0.0  0:00.00  /usr/sbin/NetworkManager --no-daemon
  1209 root      20   0  448M 10956  6968  S  0.0  0.0  0:00.09  /usr/sbin/NetworkManager --no-daemon
  1139 root      20   0  448M 10956  6968  S  0.0  0.0  0:02.34  /usr/sbin/NetworkManager --no-daemon
  1144 avahi     20   0 30220  1564  1300  S  0.0  0.0  0:00.44  avahi-daemon: running [cnf004.local]
  1148 dbus     20   0 28824  1772  1352  S  0.0  0.0  0:00.44  /bin/dbus-daemon --system --address=systemd: --nofork
  1166 root      20   0  198M  1236   776  S  0.0  0.0  0:00.00  /usr/sbin/gssproxy -D
  1167 root      20   0  198M  1236   776  S  0.0  0.0  0:00.00  /usr/sbin/gssproxy -D
  1168 root      20   0  198M  1236   776  S  0.0  0.0  0:00.00  /usr/sbin/gssproxy -D
  1169 root      20   0  198M  1236   776  S  0.0  0.0  0:00.00  /usr/sbin/gssproxy -D
  1170 root      20   0  198M  1236   776  S  0.0  0.0  0:00.00  /usr/sbin/gssproxy -D
  1164 root      20   0  198M  1236   776  S  0.0  0.0  0:00.30  /usr/sbin/gssproxy -D

F1Help F2Setup F3Search F4Filter F5Free F6SortBy F7Nice F8Nice F9Kill F10Quit
```

Slurm: Scripts sbatch

Permiten:

- Ejecutar un scripts batch sin necesidad de estar siempre conectado
- Monitorear el estado de la tarea
- Monitorear recursos
- Monitorear estado de la cola
- Monitorear estado de las particiones

Módulos

Leftraru carga sus programas mediante “módulos”

- Permite varias versiones del mismo software
- No genera conflictos entre versiones
- El software está centralizado

Para cargar un módulo:

```
module load intel
```

module avail: muestra todos los nodos disponibles

module list: lista todos los módulos cargados

module unload: quita un módulo previamente cargado

SBATH Ejemplo de script básico test.sh

Utilizar su editor por consola preferido: vim, nano

```
#!/bin/bash
#SBATCH --job-name=ejemplo
#SBATCH --partition=slims
#SBATCH -n 1
#SBATCH --output=archivo_%j.out
#SBATCH --error=archivo_%j.err
#SBATCH --mail-user=usuario@gmail.com
#SBATCH --mail-type=ALL

module load intel

sleep 10
```

Ejecución:
sbatch test.sh

SBATCH Job Array

Slurm permite enviar y administrar millones de trabajos similares de una sola vez.

- Todos los trabajos deben tener las mismas condiciones iniciales
- Provee de variables para controlar la ejecución de los jobs

```
#!/bin/bash
#SBATCH --job-name=sleep-test
#SBATCH --partition=slims
#SBATCH -n 1
#SBATCH --output=st_%j.out
#SBATCH --error=st_%j.err
#SBATCH --array=1-10
#SBATCH --mail-user=usuario@dominio.cl
#SBATCH --mail-type=ALL

./ejecutable entrada_$(SLURM_ARRAY_TASK_ID)
```


SBATCH Ejemplo OpenMP

```
#!/bin/bash
#SBATCH --job-name=ejemplo
#SBATCH --partition=slims
#SBATCH --ntasks=1 # igual a parámetro "-n"
#SBATCH --cpus-per-task 20 # "-c"
#SBATCH --output=archivo_%j.out
#SBATCH --error=archivo_%j.err
#SBATCH --mail-user=usuario@gmail.com
#SBATCH --mail-type=ALL

export OMP_NUM_THREADS=20

./ejecutable
```

Ejecución:
sbatch test.sh

SBATCH Intel MPI

```
#!/bin/bash
#SBATCH --job-name=ejemplo
#SBATCH --partition=slims
#SBATCH -n 2
#SBATCH --output=ejemplo_impi_%j.out
#SBATCH --error=ejemplo_impi_%j.err
#SBATCH --mail-user=usuario@gmail.com
#SBATCH --mail-type=ALL

module load intel impi

srun ./hola_mundo
```

```
Hello from thread 00 out of 1 from process 00 out of 2 on
cn005
Hello from thread 00 out of 1 from process 01 out of 2 on
cn005
```

Límites generales cuentas de usuario

- 120 CPUs
- 80 GB almacenamiento Lustre
- Walltime 3 días

Solicitud de cuentas:

visitar: <http://www.nlhpc.cl> (Servicios > Servicios para la Academia > Formulario)

ó solicitar información a info@nlhpc.cl

SSH: Secure Shell - Introducción

- Protocolo seguro de acceso remoto a máquinas
- El cliente también se llama SSH
- Permite:
 - Ejecutar órdenes
 - Redirigir las "X" (sistema ventanas Unix, Linux)
 - Transferir archivos mediante FTP cifrado (sftp)
 - Transferencia de archivos bidireccional (scp)
 - Túneles de conexión, entre otras cosas
- Funciona comúnmente en el puerto TCP 22
- Admite múltiples tipos de cifrado

Funcionamiento Protocolo ssh

1. Se determina identidad de cliente y servidor
2. Establecimiento de canal seguro (cifrado 256 bits)
3. Cliente inicia sesión (autenticación) en el servidor

Dos métodos de autenticación:

- Por clave: mediante credenciales (usuario y password)
- Por llave: el cliente instala su llave pública en el servidor

Instalación ssh

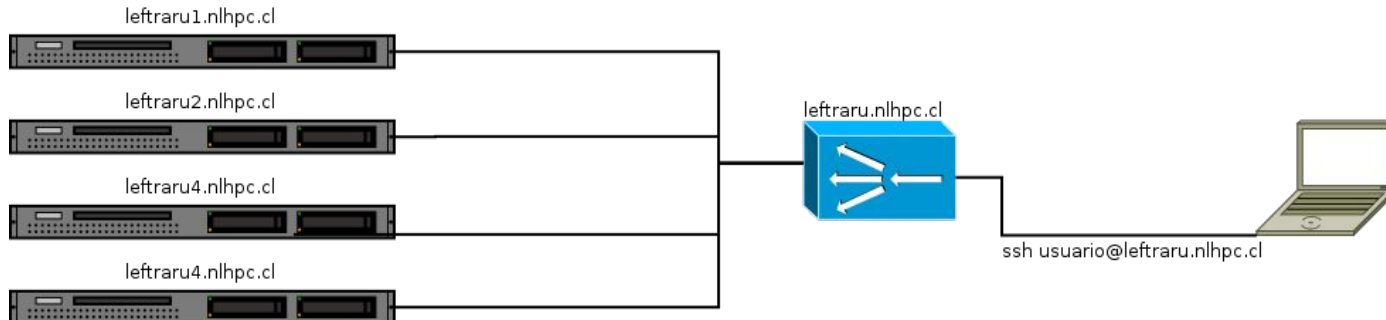
- Linux y Mac OSX: clientes integrados
- Windows:
 - Putty: cliente SSH: <http://www.putty.org>
 - WinSCP: transferencia de archivos: <http://winscp.net>

Consideraciones antes de conectar

- Balanceo de carga ssh:

El usuario siempre hará login al nodo menos utilizado (OJO: uso de screen)

- Límite de trabajos en nodos login: 10 minutos
 - Esto afecta también a las transferencias de archivos grandes



ssh: copia de archivos desde línea de comandos

ssh permite copiar desde y hacia equipos remotos

scp archivo.txt usuario@leftrarunlhpc.cl:~/ copia desde local hacia el servidor

scp usuario@leftrarunlhpc.cl:~/archivo.txt ~/ copia desde el servidor hacia carpeta local

scp -C archivo.txt usuario@leftrarunlhpc.cl:~/ “-C” utiliza compresión, envío más rápido

scp -r Directorio usuario@leftrarunlhpc.cl:~/ copia un directorio completo

Nota: ~/ ruta relativa, indica la raíz del home del usuario

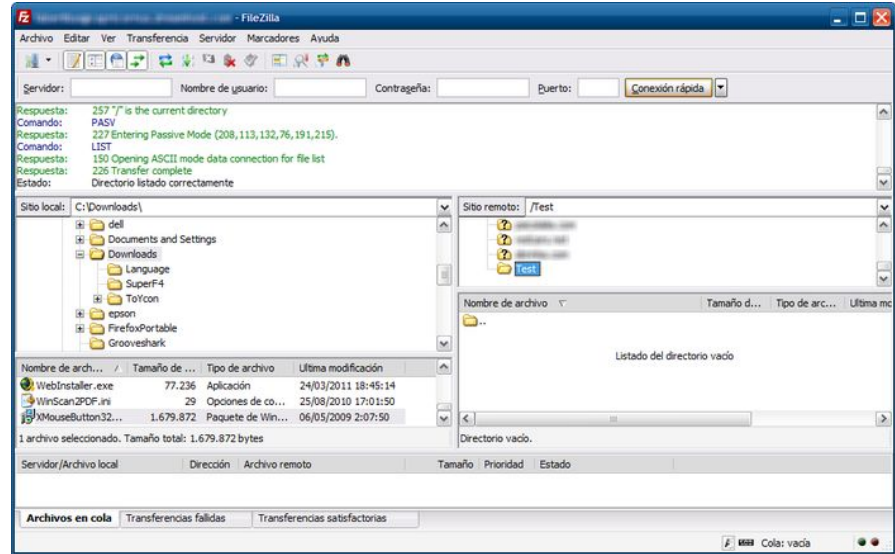
ssh copia de archivos GUI



Filezilla: <http://filezilla-project.org>

Otras opciones:

- WinSCP
- Midnight Commander (mc)



bash: screen

Cada vez que cerramos sesión, todos los procesos abiertos se cierran. Incluso si los procesos fueron puestos en background.

Screen “desacopla” una sesión completa y la envía al fondo, lo que permite volver a retomarla incluso si la conexión se perdió.

Sólo debe ejecutar screen y luego ejecutar lo que se requiera.

Para desacoplar sólo debe presionar la siguiente secuencia: **control + a** y luego **d**

Para listar las sesiones abiertas: `screen -list`

Para remontar una sesión: `screen -r id`, donde `id` es el id de la sesión

Nota: leftraru cuenta con 4 nodos de login, por lo que usted debe recordar donde ejecutó screen

Fin Introducción: Manos a la obra

¿DUDAS?

Ejercicios básicos:

Ejecute el siguiente comando: `cp -r /home/courses/ejemplos/ .`

1.- Ejecute:

- `srun --reservation=curso -N 1 hostname`
- `srun --reservation=curso -N 2 hostname`
- `srun --reservation=curso -N 3 hostname`

2.-

- Cree un script bash que ocupe 40 cores y ejecute el comando `sleep 1000`
- Lance el script
- Intente lanzarlo nuevamente ¿qué sucede?
- `scancel`: cancelar tareas

3.- Cree un script que reserve 1 nodo completo de forma exclusiva y ejecute el comando `sleep 1000`

4.- Cree un script que lance 4 trabajos, pero sólo dos por cada nodo

5.- Cree un script que lance un trabajo de un proceso en un nodo “fat”

```
Un amigo fiel: man <- sistema de ayuda en los programas linux  
man sbatch
```

Ejercicio 6: job array

En el siguiente ejercicio usted jugará con la precisión de cálculo del número pi
Ingrese el directorio ejemplos/ejercicio_6, cree el siguiente script y luego ejecútelos:

```
#!/bin/bash
#SBATCH --job-name=pi-test
#SBATCH --partition=slims
#SBATCH -n 1
#SBATCH --output=st_%j.out
#SBATCH --error=st_%j.err
#SBATCH --array=1-10
#SBATCH --mail-user=usuario@mail.cl
#SBATCH --mail-type=ALL
#SBATCH --reservation=curso

module load intel impi
PRECISION=$( echo "$SLURM_ARRAY_TASK_ID*100000000" | bc )
srun ./pi_mpi.exe $PRECISION
```

Vigile la salida: watch -n 1 squeue

Una vez terminado los procesos, analice los archivos de salida (cat *.out)

¿Nota alguna diferencia?

Ejercicio 7: Intel - MPI

- 1.- Ingrese a la carpeta ejemplos/ejercicio_7
- 2.- Compile el programa hello.c con intel mpi
mpiicc hello.c -o hello -fopenmp
Ejecútelo con srun -n 1
Ejecútelo con srun -n 20
Ejecútelo con srun -n 21
Ejecútelo con srun -c 1
Ejecútelo con srun -c 2
Ejecútelo con srun -c 20
Ejecútelo con srun -c 21 (¿qué ocurre?)
Ejecútelo con srun -n 1 -c 20
Ejecútelo con srun -n 2 -c 20

Ejercicio 8: Intel - MPI

- 1.- Ingrese a la carpeta ejemplos/ejercicio_8
- 2.- Compile el programa hello.c con intel mpi
`mpiicc hello.c -o hello -fopenmp`
- 3.- Ejecute la aplicación con 2 procesos MPI y 20 OpenMP
- 4.- Analice los archivos de salida

Ejercicio 9: SBATCH Control tareas por RAM

- Ingrese al directorio `~/ejemplos/ejercicio_9`
- Vea el contenido del ejemplo `script.sh` (`cat script.sh`)
- Ejecute y vigile la tarea
 - ¿Cuántos nodos ocupa?
 - ¿Cuanta RAM ocupa?
 - **Cancele la tarea**
- Edite nuevamente `script_slurm.sh`
 - Añada la línea: `#SBATCH --mem-per-cpu=8192`
- Vuelva a lanzar `script_slurm.sh`
 - ¿Qué ocurre?
 - ¿Cuanta RAM ocupa?
 - ¿Cuántos nodos ocupa ahora la tarea?

Ejercicio 10: Compilación - mkl

Intel MKL es una librería de optimización matemática que utiliza vectorización para mejorar el rendimiento

1.- Ingrese a la carpeta ejemplos/ejercicio_10 y liste los archivos (ls -l)

2.- Ingrese a matrix/linux

- El makefile de este ejemplo puede producir dos binarios: icc y gcc
- Para producir el binario con gcc: make gcc, producirá un archivo matrix.gcc
- Para producir el binario con icc: make gcc, producirá un archivo matrix.icc
- Ejecute ambos binarios
 - ./matrix.gcc
 - ./matrix.icc

3.- Retroceda a ingrese al directorio matrix/linux_mkl (cd ../linux_mkl)

- El makefile de este ejemplo puede generar un binario con intel MKL
- Para producir el binario con mkl: make mkl, producirá un archivo matrix.mkl
- OMP_NUM_THREADS es una variable que define la cantidad de procesos openmp
- export OMP_NUM_THREADS=1
- ./matrix.mkl

Ejercicio 10: Compilación - mkl

Ejercicio práctico:

1. Cree un script sbatch que ejecute el binario matrix.mkl, utilizando 1 proceso y 20 threads openmp
2. Modifique su script para lanzar ahora 2 procesos y 20 threads openmp

Hilos openmp: export OMP_NUM_THREADS=XX donde XX es el número de hilos

Compare los resultados

Ejercicio 11: Reservas mal hechas

- Ingrese a la carpeta ~/ejemplos/ejercicio_11
- Analice el archivo script.sh
 - ¿Cuántas CPU está reservando?
- Ejecútelo
- Vigile su tarea con ganglia
 - ¿Cuántas CPU está utilizando?
 - Cancele la tarea
- Modifique su script para que utilice las CPU que corresponden

Compilación y Optimización

- Ingrese a la carpeta “ejemplos” de su home
- Cargue el módulo intel
- Compile el ejemplo pi.c sin ningún flag de optimización y ejecútelo
- Compile el ejemplo pi.c con el flag “-O3” y ejecútelo
 - ¿Cuál fue la diferencia?
- Compile el ejemplo pi_openmpi.c
 - Ejecútelo con 1 sólo proceso
 - Ejecútelo con 10 procesos
 - Ejecútelo con 60 procesos
 - ¿Cuál es la diferencia?

FIN

¿DUDAS?

soporte@nlhpc.cl